

Ex. 1.

1.1 $y = \theta_0 + \theta_1 x + \theta_2 x^2 + \dots + \theta_p x^p + \varepsilon$ $\varepsilon \sim \mathcal{N}(0, \sigma^2)$

Dataset: $\{(x_i, y_i)\}_{i=1}^N$

(a) $\varepsilon = y - (\theta_0 + \theta_1 x + \dots + \theta_p x^p)$

$$p(\varepsilon) = \prod_{n=1}^N p(\varepsilon_n) = \left(\frac{1}{\sqrt{2\pi}\sigma} \right)^N \prod_{i=1}^N \exp \left(-\frac{\varepsilon^2}{2\sigma^2} \right)$$

$$p(y|\theta, \varepsilon) = \mathcal{N}(y; \theta, \sigma^2) = \left(\frac{1}{\sqrt{2\pi}\sigma} \right)^N \cdot \prod_{i=1}^N \exp \left(\frac{(y_i - (\theta_0 + \dots + \theta_p x_i^p))^2}{2\sigma^2} \right)$$

$$\begin{aligned} \rightarrow \mathcal{L}\{p(y|\theta)\} &= -N \log(\sqrt{2\pi}\sigma) - \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - (\theta_0 + \dots + \theta_p x_i^p))^2 \\ &= -N \log(\sqrt{2\pi}\sigma) - \frac{1}{2\sigma^2} \|y - X\theta\|^2 \end{aligned}$$

where $y \in \mathbb{R}^N$, $X \in \mathbb{R}^{N \times p}$, $\theta \in \mathbb{R}^p$

(b) Maximizing \mathcal{L} in this case is equal to find the least squares solution, since the noise is gaussian

Hence $\frac{d\|y - X\theta\|^2}{d\theta} = 0 \rightarrow \frac{d}{d\theta} (y - X\theta)^T (y - X\theta) = 0$

$$\rightarrow \frac{d}{d\theta} (y^T y - 2y^T X\theta + \theta^T X^T X\theta) = 0 \rightarrow$$

$$\rightarrow -2y^T X + 2X^T X\theta = 0 \rightarrow \boxed{\hat{\theta} = (X^T X)^{-1} X^T y}$$

$$(c) Y = \theta_0 + \theta_1 x + \theta_2 x^2 + \varepsilon, \quad \theta_0 = 0.3, \theta_1 = -0.1, \theta_2 = 0.5$$

$$\sigma^2 = 0.001$$

$$x = [-0.5, 0.2, 0.1]$$

$$\rightarrow Y = [0.4716, 0.4696, 0.3624, 0.3849, 0.2789, 0.2237, 0.2652, 0.2785]$$

$$(d) \text{ for } P=2: \quad \hat{\theta} = [0.2939, -0.1009, 0.5783] \quad L = -12,93$$

$$\text{for } P=1: \quad \hat{\theta} = [0.3103, -0.2743] \quad L = -10,12$$

$$\text{for } P=7: \quad \hat{\theta} = [0.3, 0.2, -1.5, -58.7, -19.6, 752.8, 2312] \\ L_7 = -7,3515$$

We can notice that with too many parameters ($P=7$) the weights get very big and we completely overfit. Also it is possible that the likelihood has max-value for $P=7$.

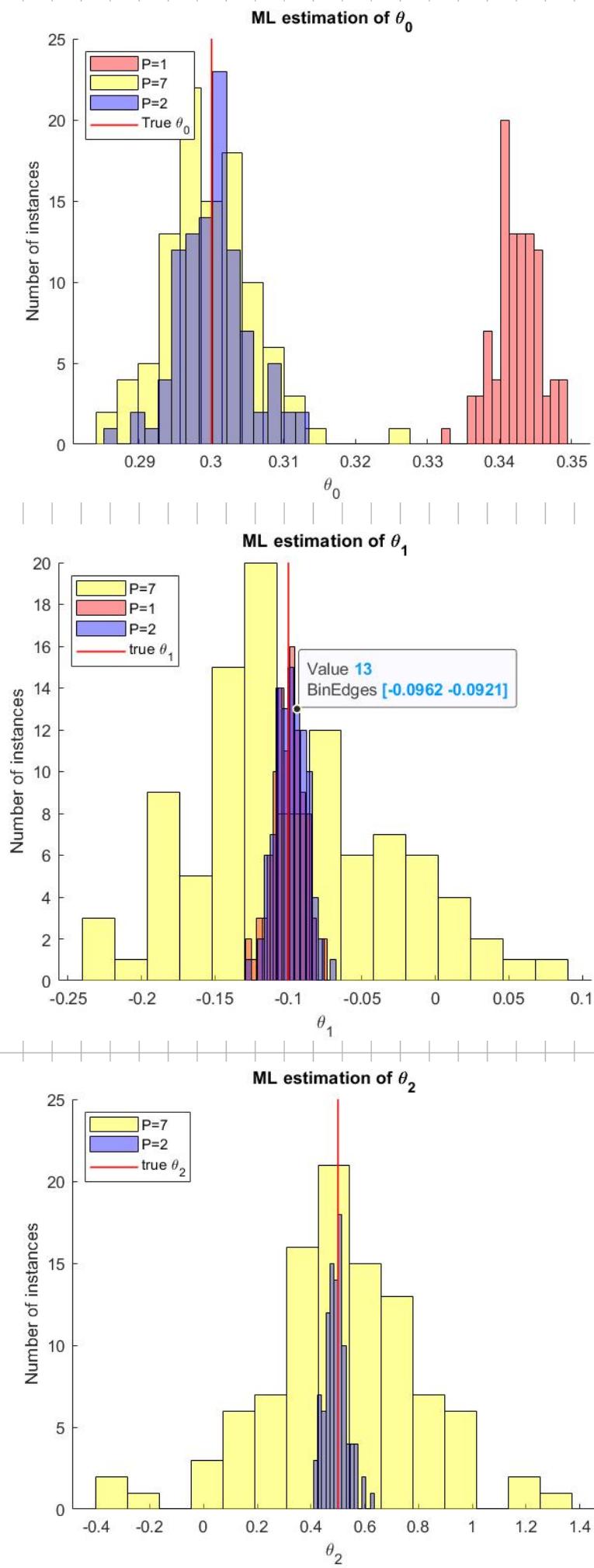
$$(e) \quad \text{For } P=7: -376.6137 \quad \text{using the parameters calculated in (d)}$$

$$P=2: -168.5593$$

$P=7: -1,0 \cdot 10^7 \rightarrow \text{enormous since the weights computed in (d) are gigantic.}$

Here we can see that the maximum likelihood value is achieved at $P=2$, as expected, since $P=1$ underfitted and $P=7$ overfitted, when giving new test data the result explains itself alone.

(f)



As expected, it is possible to see how the the number of parameters is more consistent for the case P=2. We have got the lowest variance for P=2 and the parameters are concentrated around the true value. For P=7 we have got the highest variance in the values of the estimation, this is caused by the fact that we are overfitting and and the weights of the regression get huge (at least for the thetas from 3 to 7, that are not shown here since we have already seen that they are very big, and it does not make sense to plot them since we cannot compare them with P=1 or P=2).

$$[1.2] \quad \theta = (\theta_0, \dots, \theta_p, \dots, \theta_P)^T, \quad p(\theta) = \frac{1}{\sqrt{2\pi \Sigma_0}} \exp\left(-\frac{1}{2}(\theta - \mu_0)^T \Sigma_0^{-1} (\theta - \mu_0)\right)$$

$$\Sigma_0 = I, \quad \mu_0 = 0$$

$$(a) \quad p(\theta | Y) = \frac{p(Y|\theta) p(\theta)}{p(Y)} \quad \text{d} \quad p(Y|\theta) p(\theta) = \\ = \left(\frac{1}{\sqrt{2\pi \sigma^2}} \right)^N \cdot \prod_{i=1}^N \exp\left(\frac{(Y_i - (\theta_0 + \dots + \theta_p x_i^p))^2}{2\sigma^2} \right) \cdot \frac{1}{\sqrt{2\pi} |\Sigma_0|} \exp\left(-\frac{1}{2} (\theta - \mu_0)^T \Sigma_0^{-1} (\theta - \mu_0) \right)$$

$$\mathcal{L}(p(\theta | Y)) = -N \log(\sqrt{2\pi} \sigma^2) - \sum_{i=1}^N \frac{(Y_i - (\theta_0 + \theta_1 x_i + \dots + \theta_p x_i^p))^2}{2\sigma^2} \\ - \log \sqrt{2\pi} - \frac{1}{2} \log(|\Sigma_0|) - \frac{1}{2} (\theta - \mu_0)^T \Sigma_0^{-1} (\theta - \mu_0)$$

$$(b) \quad \frac{d \mathcal{L}(p(\theta | Y))}{d\theta} = 0$$

$$\rightarrow \frac{\partial}{\partial \theta} \left(-\frac{1}{2\sigma^2} \|Y - X\theta\|^2 - \frac{1}{2} (\theta - \mu_0)^T \Sigma_0^{-1} (\theta - \mu_0) \right) = 0$$

$$\rightarrow \frac{\partial}{\partial \theta} \left(-\frac{1}{2\sigma^2} (Y - X\theta)^T (Y - X\theta) - \frac{1}{2} \theta^T \theta \right) = 0$$

$$\rightarrow \frac{\partial}{\partial \theta} \left(-\frac{1}{2\sigma^2} (Y^T Y - Y^T X\theta - \theta^T X^T Y + \theta^T X^T X\theta) - \frac{1}{2} \theta^T \theta \right) = 0$$

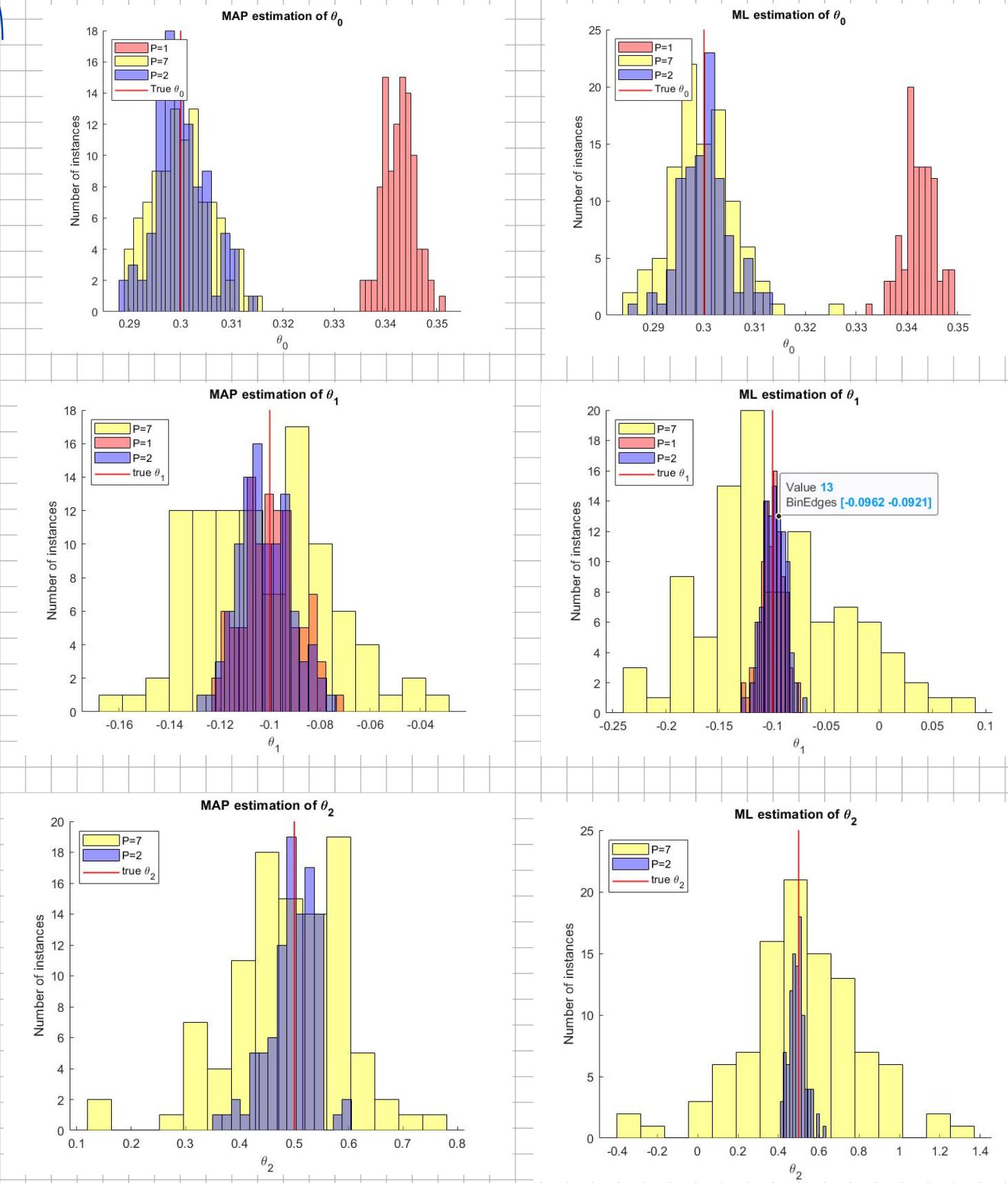
$$\frac{1}{2\sigma^2} (2X^T Y - 2X^T X\theta) - \theta = 0$$

$$X^T Y - X^T X\theta - \theta \sigma^2 = 0 \rightarrow (X^T X + \sigma^2 I)\theta = X^T Y$$

$$\rightarrow \boxed{\hat{\theta} = (X^T X + \sigma^2 I)^{-1} X^T Y}$$

(c) We can see that in the MAP case the variance of the parameters is smaller compared to the ML estimate. Not only, in fact, if we have at one instance of $\hat{\theta}$ for $P=7$: $\hat{\theta} = [0.2963, -0.1160, 0.4756, 0.1615, 0.0673, -0.6167, 0.1755, -0.2990]$, we can see that the parameters from $\theta_3 - \theta_7$ are much smaller now compared to the ML case, this is caused by the regularization effect of MAP (in this case is ridge regression) on the parameters.

(d)



It is noticeable the fact that the parameters for MAP estimation are much more concentrated around the true parameter, the scale of the X-axis is much smaller compared to the one of the ML estimation. This is caused by the σ^2 , that is proportional to the regularization factor (λ) for ridge regression. By increasing λ we strengthen the regularization, but in this case also the noise (ϵ) of the dataset.
 In conclusion MAP prevents overfitting by keeping the estimated parameters small. (If we set σ^2 to e.g 100, we can see that the parameters concentrate around 0, get very small, and of course wrong).

$$\boxed{1.3} \quad Y = X\theta + \varepsilon, \quad p(\varepsilon) \sim \mathcal{N}(0, \sigma_\varepsilon^2), \quad p(\theta) \sim \mathcal{N}(\mu_p, \sigma_p^2)$$

$$(a) \quad \varepsilon = X\theta - Y \rightarrow p(Y|\theta) = \frac{1}{\sqrt{2\pi\sigma_\varepsilon^2}} \exp\left(-\frac{(X\theta - Y)^2}{2\sigma_\varepsilon^2}\right)$$

$$\mathcal{L}(\theta) = p(Y|\theta)$$

If we have a Dataset $\{Y_i\}_{i=1}^N$

$$\rightarrow \mathcal{L}(\theta) = \left(\frac{1}{2\pi\sigma_\varepsilon^2}\right)^N \prod_{i=1}^N \exp\left(-\frac{(X\theta - Y_i)^2}{2\sigma_\varepsilon^2}\right) =$$

$$= \left(\frac{1}{2\pi\sigma_\varepsilon^2}\right)^N \exp\left(-\frac{1}{2\sigma_\varepsilon^2} \sum_{i=1}^N X^2\theta^2 - 2X\theta Y_i + Y_i^2\right) =$$

$$= \left(\frac{1}{2\pi\sigma_\varepsilon^2}\right)^N \exp\left(-\frac{1}{2\sigma_\varepsilon^2} (N\theta^2 - 2X\theta N\bar{Y} + N\bar{Y}^2)\right)$$

$$(b) \quad p(\theta|Y) = \frac{p(Y|\theta)p(\theta)}{\underbrace{p(Y)}_{\varepsilon p(Y)}} \propto p(Y|\theta)p(\theta)$$

$$\rightarrow p(\theta|Y) \propto \left(\frac{1}{\sqrt{2\pi\sigma_\varepsilon^2}}\right)^N \exp\left(-\frac{1}{2\sigma_\varepsilon^2} (N\theta^2 - 2X\theta N\bar{Y} + N\bar{Y}^2)\right).$$

$$\cdot \frac{1}{\sqrt{2\pi\sigma_p^2}} \exp\left(-\frac{(\theta - \mu_p)^2}{2\sigma_p^2}\right) = \frac{1}{(2\pi\sigma_\varepsilon^2)^{N/2} (2\pi\sigma_p^2)^{1/2}}.$$

$$\cdot \exp\left(-\frac{1}{2\sigma_\varepsilon^2} (N\theta^2 - 2X\theta N\bar{Y} + N\bar{Y}^2) - \frac{\theta^2 - 2\theta\mu_p + \mu_p^2}{2\sigma_p^2}\right)$$

$$\rightarrow \log p(\theta|Y) \propto -\frac{N}{2} \log(2\pi\sigma_\varepsilon^2) - \frac{1}{2} \log 2\pi\sigma_p^2 +$$

$$-\frac{1}{2\sigma_\varepsilon^2} (N\theta^2 - 2X\theta N\bar{Y} + N\bar{Y}^2) - \frac{\theta^2 - 2\theta\mu_p + \mu_p^2}{2\sigma_p^2}.$$

$$(c) \log N(\theta | M, \sigma^2) = -\frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (\theta^2 - 2M\theta + M^2)$$

Take the "exp" part $\rightarrow -\frac{1}{2\sigma_{\epsilon}^2\sigma_p^2} (\theta^2(Nx^2\sigma_p^2 + \sigma_{\epsilon}^2) +$

$$-2\theta(xN\bar{y}\sigma_p^2 + M_p\sigma_{\epsilon}^2) + (N\bar{y}\sigma_p^2 + M_p^2\sigma_{\epsilon}^2))$$

The expression in (b) is still a Gaussian, with mean M, σ_{ϵ}^2 .

Since the conjugate of a Gaussian is a Gaussian, if we multiply them we still get a Gaussian distribution with different mean and variances.

Like the Beta distribution is conjugate to the Binomial distribution.

$$(d) -\frac{N}{2} \log(2\pi\sigma_{\epsilon}^2) - \frac{1}{2} \log(2\pi\sigma_p^2) - \underbrace{\frac{1}{2\sigma_{\epsilon}^2\sigma_p^2} (\theta^2(Nx^2\sigma_p^2 + \sigma_{\epsilon}^2))}_{-2\theta(xN\bar{y}\sigma_p^2 + M_p\sigma_{\epsilon}^2) + (N\bar{y}\sigma_p^2 + M_p^2\sigma_{\epsilon}^2)} + \quad \textcircled{1}$$

$$\textcircled{1} \frac{-(Nx^2\sigma_p^2 + \sigma_{\epsilon}^2)}{2\sigma_{\epsilon}^2\sigma_p^2} \left(\theta^2 + \frac{2\theta(xN\bar{y}\sigma_p^2 + M_p\sigma_{\epsilon}^2)}{N\bar{y}\sigma_p^2 + \sigma_{\epsilon}^2} \right) + \mathcal{N}$$

- . . . Sorry, too long.